# CUNY SPS DATA 621 - CTG5 - Final

*Betsy Rosalen, Gabrielle Bartomeo, Jeremy O'Brien, Lidiia Tronina, Rose Jones*

*May 25th, 2019*

## Contents

# 1  PROJECT DESCRIPTION AND BACKGROUND

## 1.1  Background

With nearly 18MM deaths in 2015, cardiovascular diseases (CVD) are the leading cause of death globally and growing in the developing world. CVD is a disease class which includes heart attacks, strokes, heart failure, coronary artery disease, arrhythmia, venous thrombosis, and other conditions. About half of all Americans (47%) have at least one of three key risk factors for heart disease: high blood pressure, high cholesterol, and smoking.

Researchers estimate that up to 90% of heart disease deaths could be prevented. Typical means of detection include electrocardiograms (ECG's), stress tests, and cardiac angiograms, all of which are expensive. Risk evaluation screenings require blood samples, which are assessed alongside risk factors like tobacco use, diet, sleep disorders, physical inactivity, air pollution, and others.

More efficient, scalable, and non-invasive means of early detection could be used to trigger medical interventions, prompt preventive care by physicians, and/or engender behavioral change on the part of those prone to or suffering from CVD. Applying data mining techniques to CVD datasets to predict risk based on existing or easy-to-collect health data could improve healthcare outcomes and mortality rates.

The Cleveland dataset is the most complete CVD dataset and is the most frequently used in data science experimentation. It has a small number of cases (n = 303) but numerous variables including the 14 most commonly selected for analysis which were are the ones selected for this analysis as well (m = 14, including the target class). Small numbers of observations are common with health data given the costs of collecting experimental data and the privacy risk considerations of observational data.

## 1.2  Research Question and Approach

A multitude of approaches and methodologies have been attempted by researchers with the aim of predicting the presence of heart disease (the target class) based on the 13 other variables most commonly selected from the Cleveland dataset.

For this project, we evaluate different classification techniques - including logistic regression, random forests, Naive Bayes, and Support Vector Machines models - that preceding researchers have found fruitful, harness synthesized data, and attempt to improve upon the models' performance.

# 2 DATA PREPARATION

## 2.1 Original Data Description

The original data set has 13 predictor variables 8 of which are categorical and range from 2 to 5 levels and the other 5 are numeric. The target variable is a binary categorical variable that indicates whether or not the patient has heart disease with 1 indicating presence of heart disease and 0 indicating no heart disease. Descriptions of each of the variables are provided in Table 1: Data Dictionary.

Table 1: Data Dictionary

| VARIABLE | DEFINITION | TYPE |
|----------|-----------|------|
| age | Age | continuous numerical predictor |
| sex | Sex. Female = 0, Male = 1 | categorical predictor |
| cp | Chest pain type. Scale of 0 to 4 | categorical predictor |
| trestbps | Diastolic blood pressure in mmHg | continuous numerical predictor |
| chol | Serum cholesterol (mg/dl) | continuous numerical predictor |
| fbs | Fasting blood sugar. Greater than 120mg/dl, value of 0 or 1 | categorical predictor |
| restecg | Resting ECG. Value of 0, 1, or 2 | categorical predictor |
| thalach | Maximum heartrate achieved from thallium test | continuous numerical predictor |
| exang | Exercise-induced angina. Value of 0 or 1 | categorical predictor |
| oldpeak | Old-peak.ST depression induced by exercise relative to rest | continuous numerical predictor |
| slope | Slope of peak exercise ST segment, value of 1, 2, or 3 | categorical predictor |
| ca | Number of major vessels (0-3) colored by fluoroscopy | categorical predictor |
| thal | Exercise thallium scintigraphic defects | categorical predictor |
| target | Response Variable - No Heart Disease = 0, Heart Disease = 1 | categorical predictor |

### 2.1.1 Original Data Summary Statistics

Table 2: Summary statistics for numeric variables in the original data set

|  | n | min | mean | median | max | sd |
|---|---|-----|------|--------|-----|-----|
| age | 303 | 29 | 54.366337 | 55.0 | 77.0 | 9.082101 |
| trestbps | 303 | 94 | 131.623762 | 130.0 | 200.0 | 17.538143 |
| chol | 303 | 126 | 246.264026 | 240.0 | 564.0 | 51.830751 |
| thalach | 303 | 71 | 149.646865 | 153.0 | 202.0 | 22.905161 |
| oldpeak | 303 | 0 | 1.039604 | 0.8 | 6.2 | 1.161075 |

Table 3: Summary statistics for categorical variables in the original dataset

| sex | cp | ca | exang | fbs | restecg | slope | target | thal |
|-----|-----|-----|-------|-----|---------|-------|--------|------|
| 0: 96 | 0:143 | 0:175 | 0:204 | 0:258 | 0:147 | 0: 21 | 0:138 | 0: 2 |
| 1:207 | 1: 50 | 1: 65 | 1: 99 | 1: 45 | 1:152 | 1:140 | 1:165 | 1: 18 |
| NA | 2: 87 | 2: 38 | NA | NA | 2: 4 | 2:142 | NA | 2:166 |
| NA | 3: 23 | 3: 20 | NA | NA | NA | NA | NA | 3:117 |
| NA | NA | 4: 5 | NA | NA | NA | NA | NA | NA |

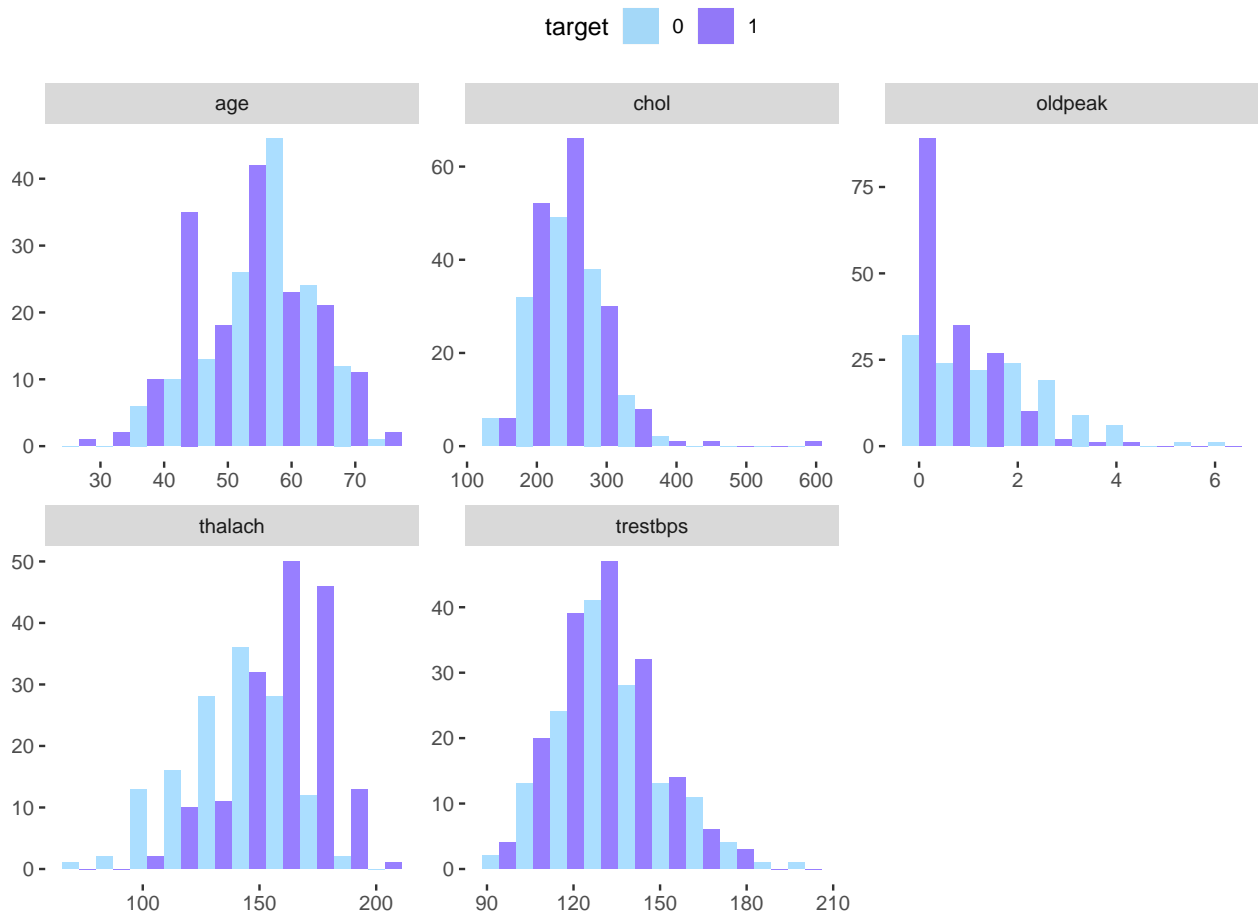## 2.1.2 Original Data Summary Statistics Graphs



Figure 1: Numeric Data Distributions as a Function of TARGET

Figure 2: Categorical Data Distributions as a Function of TARGET
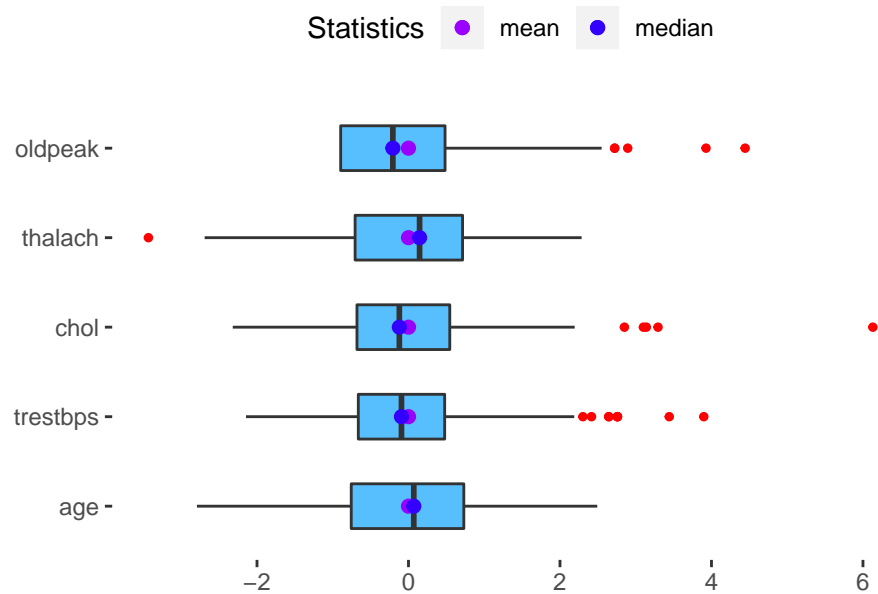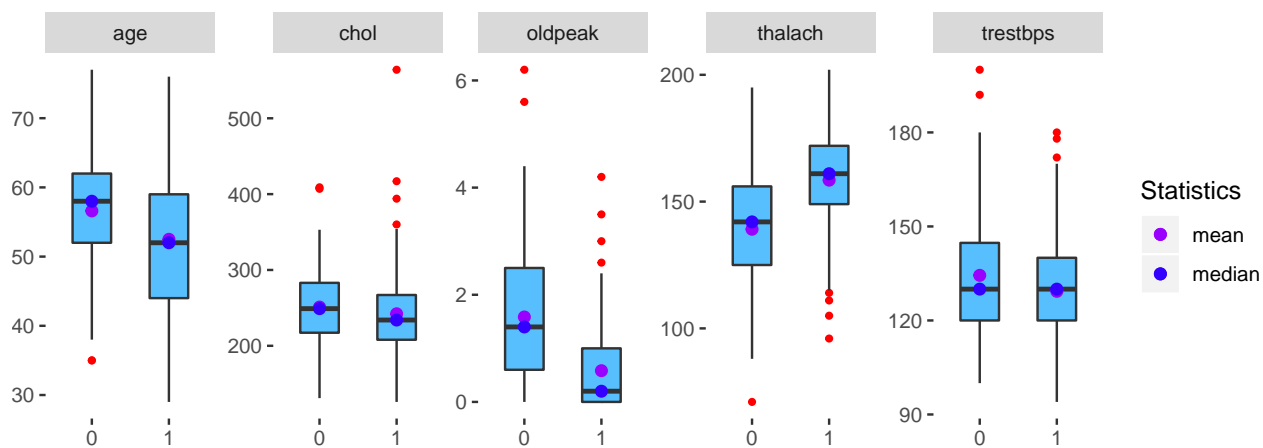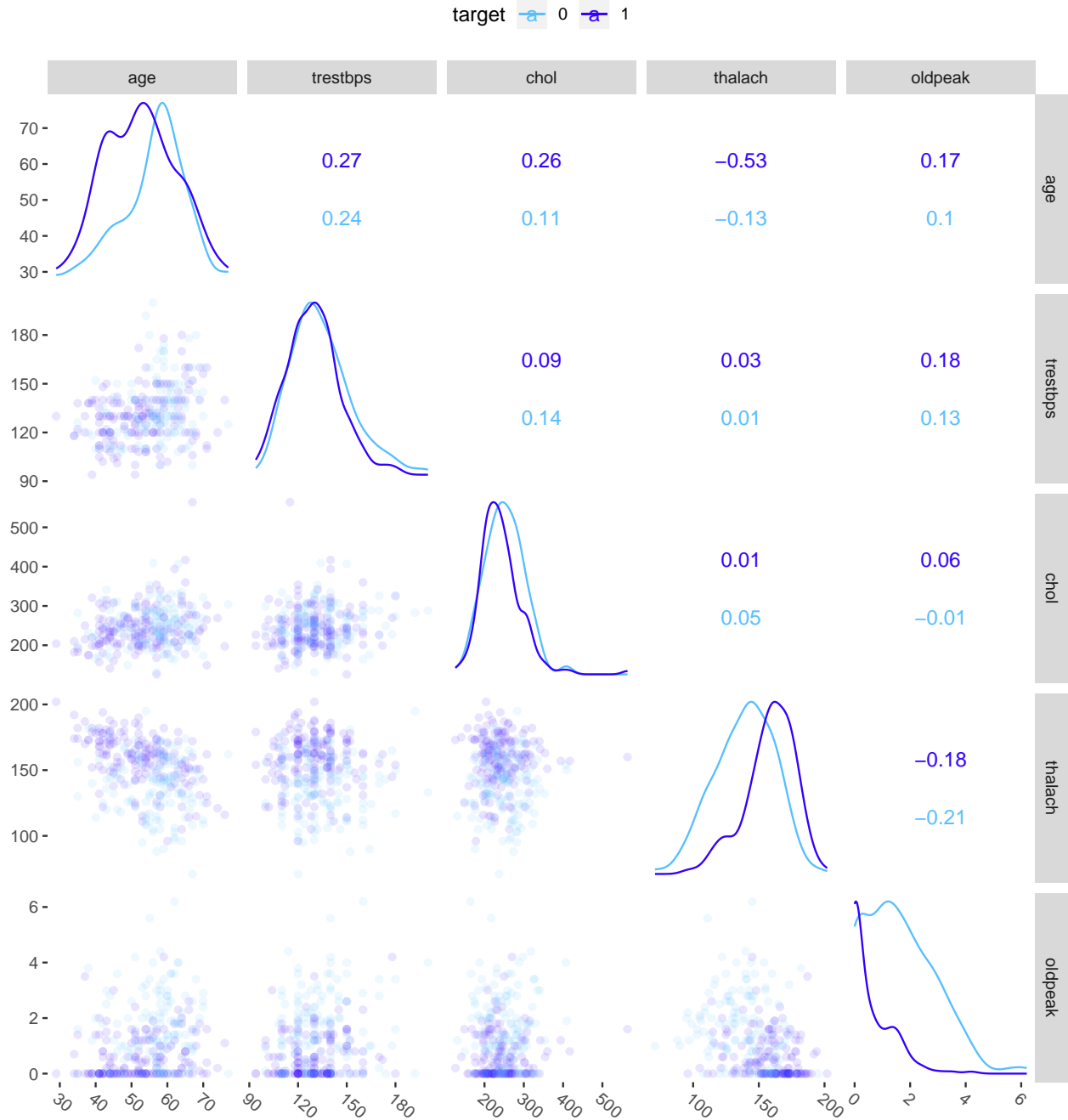
Figure 3: Scaled Boxplots for Numeric Variables



Figure 4: Linear relationship between each numeric predictor and the target

## 2.2 Cross-validation

Cross-validation is a resampling procedure to evaluate machine learning models. This approach involves randomly dividing the in-sample data into k groups or folds, of equal size. The first fold is treated as a validation set and the method is fit on the remaining k-1 folds [22]. To use nested cross-validation, the outer cross-validation provides performance assessment obtainable using a method of constructing a model, including feature selection. The inner cross-validation is used to select the features independently in each fold of the outer cross-validation [23]. We retain the evaluation score and discard the model to summarize the skill of the model using the evaluation scores.

One of the most popular cross validation techniques is Grid Search. Grid search experiments are common

in the literature of empirical machine learning, where they are used to optimize the hyper-parameters of learning algorithms. It is common to perform multi-stage automated grid experiment, however, fine resolution for optimization would be computationally expensive. Grid search experiments allocates many trials to the exploration of dimensions that may not matter and suffer from poor coverage in dimensions that are important.

On the other hand, Random search found better models in most cases and required less computational time. It is easier to carry out and more practical in terms of statistical independence of every trial. The experiment can be stopped at any time and can be added without adjustment of grid or committing larger experiment as every trial can be carried out asynchronously.

To minimize the computation time, we implement Coarse to Fine strategy. During the Coarse search phase, we use the random search technique, then, simply filter the under-performing parameter values out to run Finer search to find the best values for parameters.

## 2.3 Bootstraping 'synthetic' data

The idea of synthesizing data is to replace some or all observed values so that the statistical features of the original data are preserved. This approach can be used to anonymize data subjects, keep actual observations confidential, or comply with legal or regulatory requirements regarding identifiable information while still performing data modeling or other data-related tasks.

For our purposes, the motivation to synthesize data is to augment the size of the dataset, simulating a larger number of cases based on the original distributions of the observed data using the synthpop package in R. In this way, we can scale the data from hundreds to thousands or tens of thousands of cases with values for each variable. This enables a wider range of potential modeling techniques, and we are curious if it will allow models that can benefit from larger samples to achieve better or more stable performance.

### 2.3.1 Synthesis diagnostics

The original Cleveland dataset contains n = 303 observations. The synthesized dataset is simulated based on the same probability distribution, but is 20 times larger at n= 6,060.
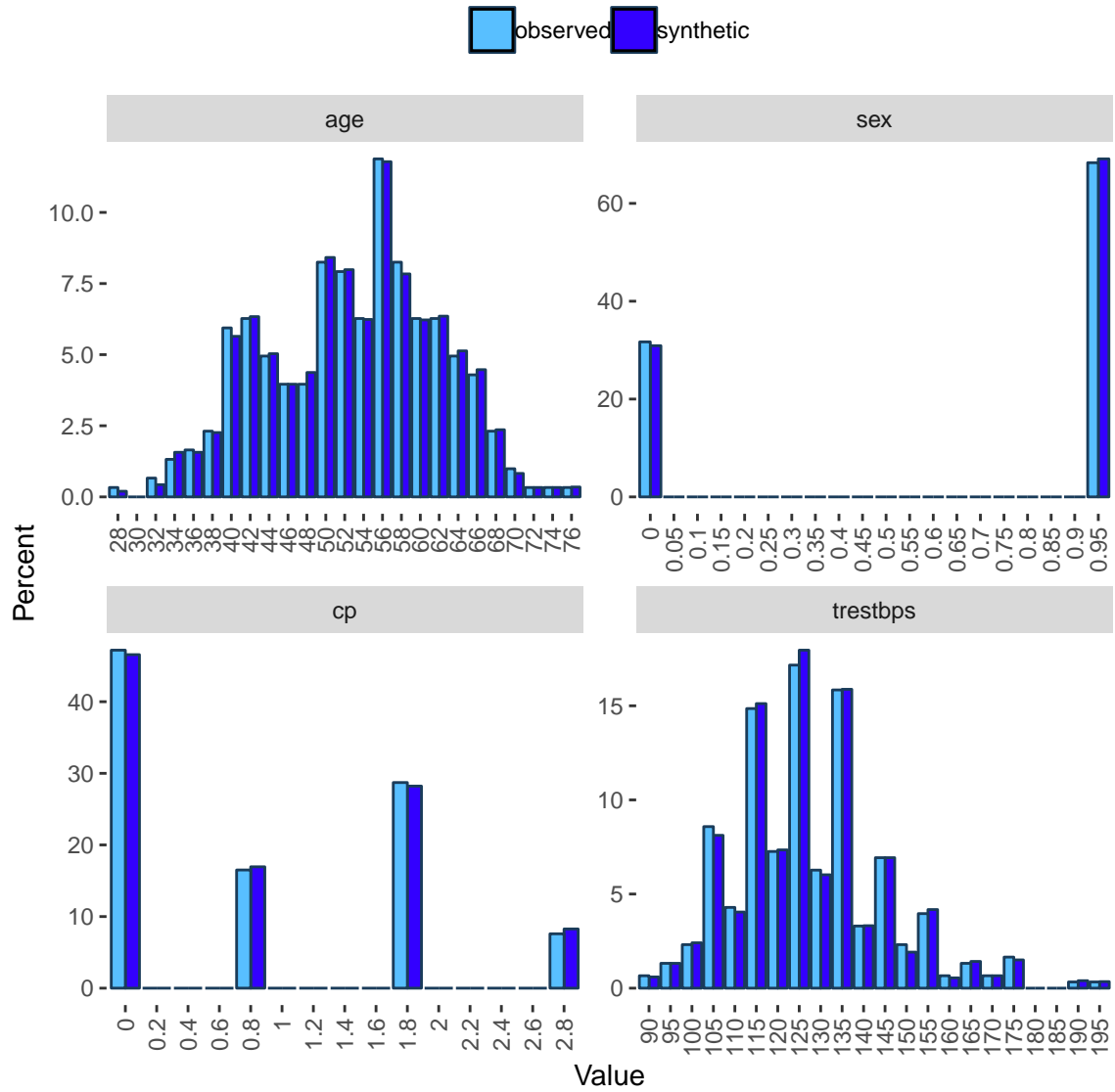
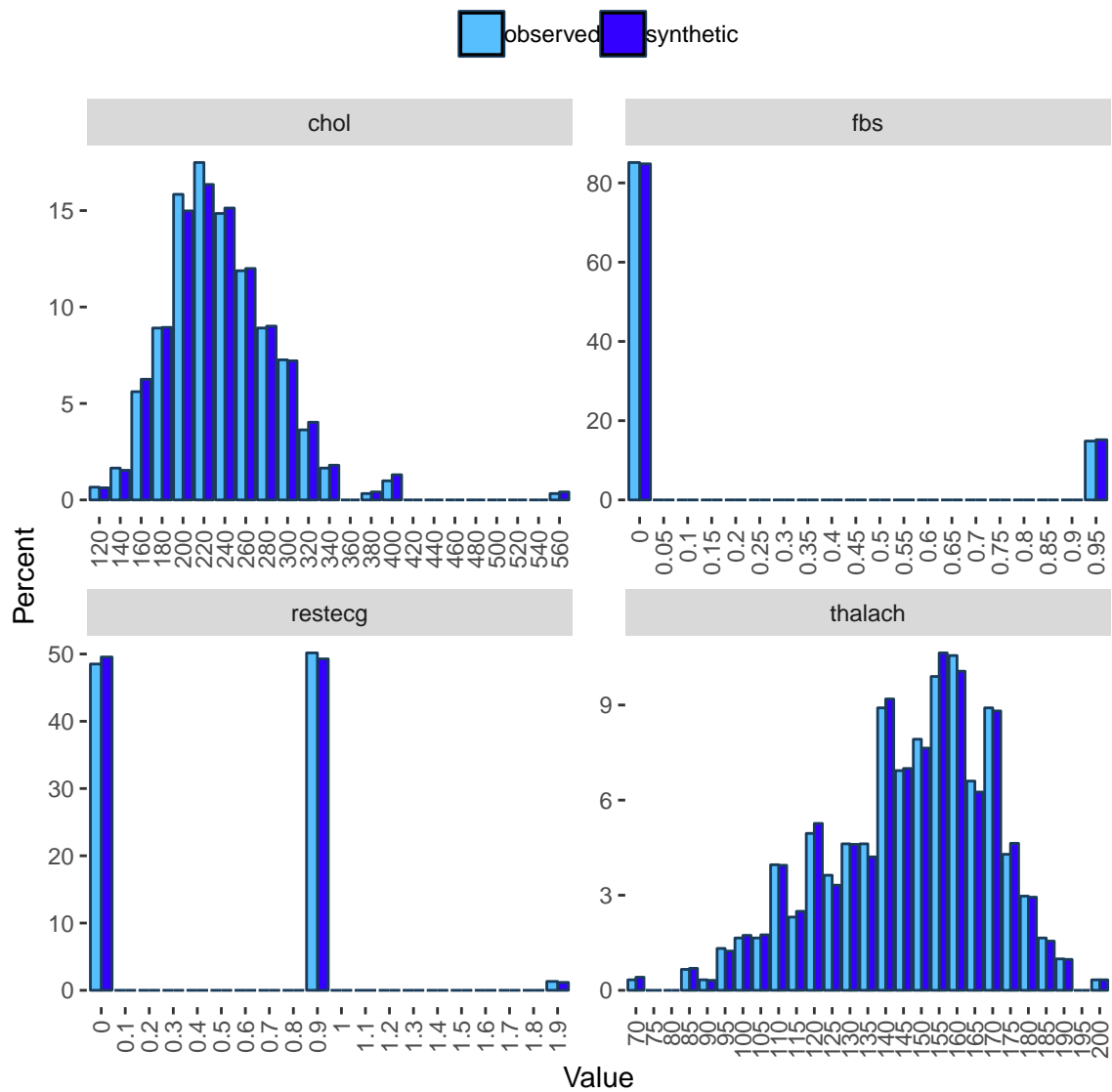Figure 5: Data Distribution - Original vs. Synthesized

Figure 6: Data Distribution - Original vs. Synthesized

Figure 7: Data Distribution - Original vs. Synthesized

Figure 8: Data Distribution - Original vs. Synthesized

### 2.3.2 Synthesized data summary statistics

Table 4: Summary statistics for numerical variables in the synthesized dataset

|          | n    | min | mean       | median | max   | sd        |
|----------|------|-----|------------|--------|-------|-----------|
| age      | 6060 | 29  | 54.432673  | 55.0   | 77.0  | 9.000207  |
| trestbps | 6060 | 94  | 131.672607 | 130.0  | 200.0 | 17.432143 |
| chol     | 6060 | 126 | 247.733828 | 243.0  | 564.0 | 53.953963 |
| thalach  | 6060 | 71  | 149.460561 | 153.0  | 202.0 | 23.045903 |
| oldpeak  | 6060 | 0   | 1.047376   | 0.8    | 6.2   | 1.147023  |

Table 5: Summary statistics for categorical variables in the synthesized dataset

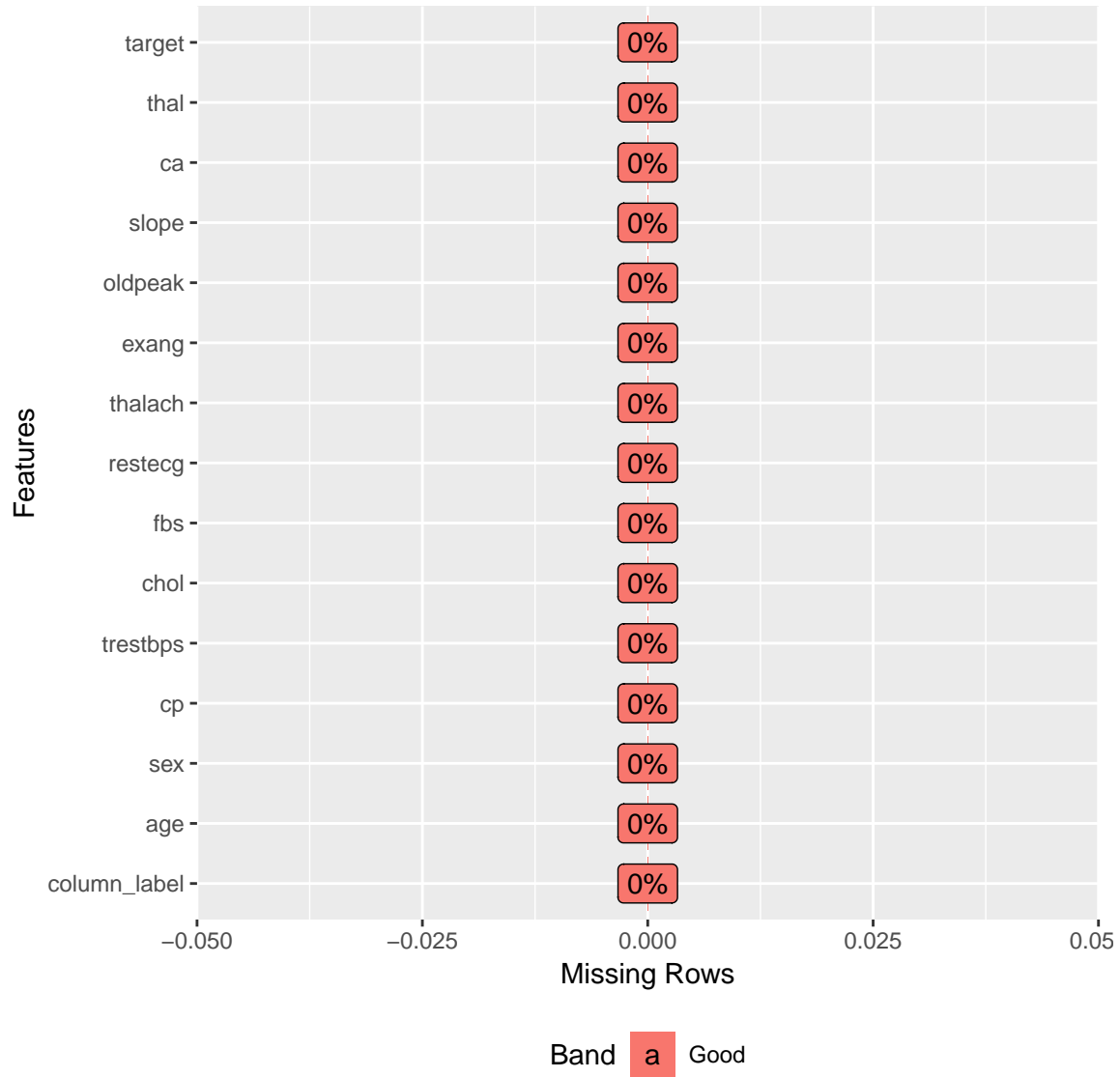| sex | cp | ca | exang | fbs | restecg | slope | target | thal |
|---|---|---|---|---|---|---|---|---|
| 0:1873 | 0:2822 | 0:3473 | 0:4039 | 0:5140 | 0:3003 | 0: 510 | 0:2828 | 0: 39 |
| 1:4187 | 1:1027 | 1:1290 | 1:2021 | 1: 920 | 1:2986 | 1:2834 | 1:3232 | 1: 423 |
| NA | 2:1710 | 2: 779 | NA | NA | 2: 71 | 2:2716 | NA | 2:3334 |
| NA | 3: 501 | 3: 410 | NA | NA | NA | NA | NA | 3:2264 |
| NA | NA | 4: 108 | NA | NA | NA | NA | NA | NA |



Figure 9: Missing data in the synthesised dataset

13

# 3  BUILDING MODELS

Our literature review highlighted a wide array of approaches to classification taken to predict the presence of heart disease using the 13 variables most commonly selected from the Cleveland dataset.

Shouman et al. compiled an exhaustive review of over 60 papers published between 2000 and 2016 that detail different classification modeling approaches built on heart disease datasets, include the Cleveland dataset:

Summarizing model performance (median accuracy) by type:

We focused on a few pieces of research in particular:

- The aim of Shouman et al.'s work is to evaluate a potential low-cost heart disease expert system risk evaluation tool leveraging non-invasive data attributes. This is explored by evaluating the Cleveland dataset alongside another dataset from Canberra not available via the UCI machine learning repository. When constrained to Cleveland's non-invasive data attributes, the best performance is seen with a combination of age, sex, and resting blood pressure. This line of research also explores integrating K-means clustering with decision tree models to improve accuracy.
- Assari et al.'s broader data-mining focus finds that SVM and Naive Bayes outperform KNN (of K=7) and Decision Tree in terms of accuracy when using 10-fold cross-validation. Its results identify the most important features as chest pain type, exercise thallium, and coronary artery disease.
- Sabay et al. seek to assess the application of ML techniques requiring more observations to the Cleveland dataset so improve its generalizability. To that end, a surrogate synthetic dataset is bootstrapped using the Synthpop package in R. Logistic Regression is found to be more accurate and stable than Random Forest and Decision Tree methods, both for the original dataset as well as a 50,000-observation surrogate. An ANN perceptron model built on a 60,000-observation surrogate dataset achieves accuracy and recall above 95%.

We selected four modeling approaches based on our review of the literature and findings:

- Logistic Regression
- Random Forest
- Support Vector Machines
- Naive Bayes

## 3.1  Logistic regression

The logistic regression model involves picking one or more variables in comparison to the target. All of the available variables were initially compared to the target variable in both the original data and the synthetic data. The result of this comparison in the original data was that the patient's sex, chest pain type, and the number of vessels colored by fluoroscopy were the most influential of the variables. In the synthetic data, the major influencers were the patient's age, sex, chest pain type, serum cholesterol, and the number of vessels colored by fluoroscopy.

Multiple models were tested for their validity for each type of data. For the original data, the first model tested was looking at the individual influence of patient's sex, chest pain type, and the number of vessels colored by fluoroscopy on the target variable. The second model tested examined the relationship between the patient's sex and chest pain type, and how both of these variables interacted with each other, and then the number of vessels colored by fluoroscopy when compared with the target variable. The third model looked at all numeric variables compared to the target, and the last model looked at all categorical (factorized) variables compared to the target.

|          | R2   | Adj..R2 | AIC      | BIC      |
|----------|------|---------|----------|----------|
| Select 1 | 0.43 | 0.41    | 235.6967 | 270.6273 |
| Select 2 | 0.44 | 0.41    | 237.9943 | 283.4041 |
| Numeric  | 0.29 | 0.27    | 283.0819 | 307.5334 |
| Factors  | 0.57 | 0.54    | 182.3636 | 248.7317 |

Across the original data, the model with the best preliminary performance was the factorized variables.

A similar process occurred with the synthetic data that resulted in three models instead of four. The first model compared the patient's age, sex, chest pain type, serum cholesterol, and the number of vessels colored by fluoroscopy to the target variable. The second model looked at all numeric variables compared to the target, and the last model looked at all categorical (factorized) variables compared to the target.

|  | R2 | Adj..R2 | AIC | BIC |
|---|---|---|---|---|
| Select | 0.40 | 0.39 | 4605.077 | 4682.913 |
| Numeric | 0.18 | 0.18 | 6090.317 | 6135.722 |
| Factors | 0.42 | 0.42 | 4417.274 | 4540.515 |

The synthetic model proves once again the factorized variables are more influential than the others tested.

To further single out the best model, two kinds of predictions were made for each model: predictions against the test data for the data the model was created from, and predictions against the test data for the data the model was not created from.

| | Accuracy | Kappa | AccuracyLower | AccuracyUpper | AccuracyNull | AccuracyPValue | Mcnemar |
|---|---|---|---|---|---|---|---|
| factors_o | 0.8266667 | 0.6494067 | 0.7218513 | 0.9043494 | 0.5466667 | 0.0000003 | 1.0 |
| factors_sXo | 0.8266667 | 0.6458409 | 0.7218513 | 0.9043494 | 0.5466667 | 0.0000003 | 0.2 |
| factors_s | 0.8158416 | 0.6277022 | 0.7953878 | 0.8350630 | 0.5333333 | 0.0000000 | 0.0 |
| factors_oXs | 0.7788779 | 0.5539281 | 0.7571204 | 0.7995487 | 0.5333333 | 0.0000000 | 0.0 |
| select_s | 0.7742574 | 0.5452934 | 0.7523548 | 0.7950915 | 0.5333333 | 0.0000000 | 0.1 |
| numeric_o | 0.7600000 | 0.5058565 | 0.6474760 | 0.8511308 | 0.5466667 | 0.0001097 | 0.0 |
| numeric_sXo | 0.7466667 | 0.4850018 | 0.6329944 | 0.8400702 | 0.5466667 | 0.0002843 | 0.6 |
| select1_oXs | 0.7458746 | 0.4914135 | 0.7231582 | 0.7676338 | 0.5333333 | 0.0000000 | 0.0 |
| select2_oXs | 0.7372937 | 0.4738234 | 0.7143561 | 0.7593078 | 0.5333333 | 0.0000000 | 0.0 |
| select_sXo | 0.7333333 | 0.4537509 | 0.6186273 | 0.8288895 | 0.5466667 | 0.0006897 | 0.2 |
| select1_o | 0.7066667 | 0.4051911 | 0.5902167 | 0.8061907 | 0.5466667 | 0.0033625 | 0.8 |
| select2_o | 0.7066667 | 0.4051911 | 0.5902167 | 0.8061907 | 0.5466667 | 0.0033625 | 0.8 |
| numeric_s | 0.6924092 | 0.3780032 | 0.6684840 | 0.7155863 | 0.5333333 | 0.0000000 | 0.0 |
| numeric_oXs | 0.6917492 | 0.3740506 | 0.6678114 | 0.7149413 | 0.5333333 | 0.0000000 | 0.0 |

The most accurate model for the original data and the synthetic data was the factorized model. The most accurate model for data not matching what the model was built on was also the factorized model. The best of the best was definitely the factorized model for the original data.

|  | 0 | 1 |
|---|---|---|
| 0 | 27 | 6 |
| 1 | 7 | 35 |

Figure 10: Confusion Matrix for Factorized Model on Original Data

## 3.2 Random forest

Random Forest consists of numerous decision trees that are generated based on bootstrap sampling from the in-sample data. Subsampling reduces the variance of trees substantially and the random feature selection decorrelates them to improve the predictive accuracy and control over-fitting.

The bootstrap resampling of the data for training each tree increases the diversity between the trees. Each tree is composed of a root node, branch nodes, and leaf nodes. For each node of a tree, the optimal node splitting feature is selected from a set of features that are again randomly selected. The final output is an ensemble of random forest trees, so that classification can be performed via majority vote.

Tuning hyperparameter values in the model significantly impact the accuracy of model performance. The three parameters that we tune in this experiment are the number of trees, maximum depth, and the number of features to randomly select. We use a coarse to fine search strategy for hyperparameter tuning. Our model used normally distributed random values for parameters for a few tries. Then, we incorporate cross-validation during the training process to evaluate the results. Upon completion, the parameters from the first pass models are used to create a refined range for parameter selection.

### 3.2.1 Hyper parameter tuning

#### 3.2.1.1 Tune using caret

The caret package in R provides an excellent facility to tune machine learning algorithm parameters. Not all machine learning algorithms are available in caret for tuning. The choice of parameters is left to the developers of the package. Only those algorithm parameters that have a large effect are available for tuning in caret. As such, only `mtry` parameter is available in caret for tuning. The reason is its effect on the final accuracy and that it must be found empirically for a dataset. The `ntree` parameter is different in that it can be as large as you like, and continues to increases the accuracy up to some point. It is less difficult or critical to tune and could be limited more by compute time available more than anything.

##### 3.2.1.1.1 Random Search

One search strategy that we can use is to try random values within a range. This can be good if we are unsure of what the value might be and we want to overcome any biases we may have for setting the parameter (like the suggested equation above). Let us try a random search for `mtry` using caret: We can see that the most accurate value for mtry was 2 with an accuracy of 0.8084378

##### 3.2.1.1.2 Grid Search

Grid search experiments are common in the literature of empirical machine learning, where they are used to optimize the hyper-parameters of learning algorithms. It is common to perform multi-stage automated grid experiment, however, fine resolution for optimization would be computationally expensive. Grid search experiments allocates many trials to the exploration of dimensions that may not matter and suffer from poor coverage in dimensions that are important. We can see that the most accurate value for `mtry` was 1 with accuracy of 0.8158102

#### 3.2.1.2 Tune Using Algorithm Tools

Some algorithms provide tools for tuning the parameters of the algorithm. For example, the random forest algorithm implementation in the randomForest package provides the tuneRF() function that searches for optimal mtry values given your data. We can see that the most accurate value for mtry was 2 with an OOBError of 0.1650165 This does not really match up with what we saw in the caret repeated cross validation experiment above, Nevertheless, it is an alternate way to tune the algorithm.

#### 3.2.1.3 Craft your own parameter search

##### 3.2.1.3.1 Tune Manually

We want to keep using caret because it provides a direct point of comparison to our previous models and because of the repeated cross validation test harness that we like as it reduces the severity of overfitting. One approach is to create many caret models for our algorithm and pass in a different parameters directly to the algorithm manually. Let's look at an example doing this to evaluate different values for `ntree` while holding mtry constant.
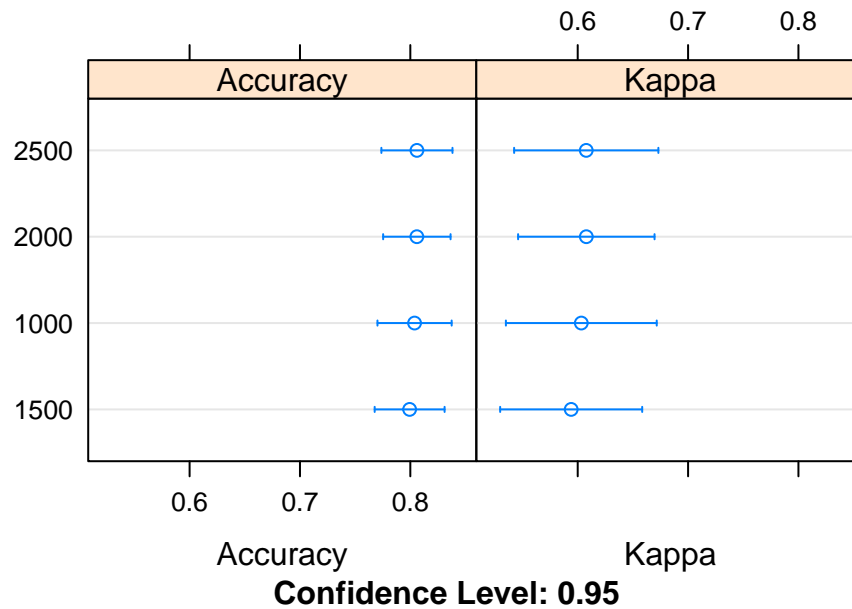
Figure 11: Random Forest Hyperparameter Tuning Manual Search

We can see that the most accuracy value for ntree was perhaps 1000 with a mean accuracy of 80.38% (a lift over our very first experiment using the default mtry value). The results perhaps suggest an optimal value for ntree between 2000 and 2500. Also note, we held mtry constant at the default value. We could repeat the experiment with a possible better mtry=2 from the experiment above, or try combinations of of ntree and mtry in case they have interaction effects.

#### 3.2.1.3.2 Extend Caret

Another approach is to create a "new" algorithm for caret to support. This is the same random forest algorithm you are using, only modified so that it supports multiple tuning of multiple parameters. A risk with this approach is that the caret native support for the algorithm has additional or fancy code wrapping it that subtly but importantly changes it's behavior. You many need to repeat prior experiments with your custom algorithm support. We can define our own algorithm to use in caret by defining a list that contains a number of custom named elements that the caret package looks for, such as how to fit and how to predict. See below for a definition of a custom random forest algorithm for use with caret that takes both an mtry and ntree parameters. Now, let's make use of this custom list in our call to the caret train function, and try tuning different values for ntree and mtry.
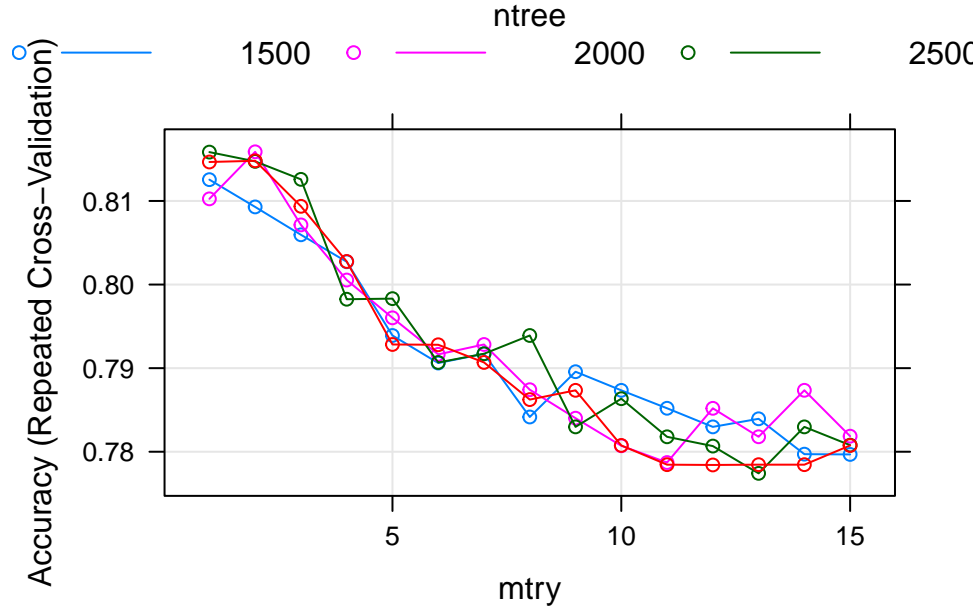
Figure 12: Random Forest Hyperparameter Tuning Custom Search

You can see that the most accurate values for ntree and mtry were 1500 and 2 with an accuracy of 84.43%. We do perhaps see some interaction effects between the number of trees and the value of ntree.

### 3.2.2 Random Forest Baseline Model

we will stick to tuning two parameters, the `mtry` and the `ntree` that have the most influence on accuracy of RF model. `mtry` is the number of variables randomly sampled as candidates at each split. `ntree` is the number of trees to grow. This should not be set to too small a number, to ensure that every input row gets predicted at least a few times. Let us create a baseline for comparison, using recommended default values for parameters: `mtry` floor(sqrt(ncol(x))) and `ntree` 500.

The baseline model used total 303 samples which includes 13 predictor variables to classify two classes in target. Resampling was done using cross-validation for 10 fold, repeating 3 times. The estimated accuracy is 0.7962489

### 3.2.3 Random Forest Final Model And Evaluation

The final model is developed by using generated synthetic data with the help of minority class data. Simply put, it takes the minority class data points and creates new data points which lie between any two nearest data points joined by a straight line. In order to do this, the algorithm calculates the distance between two data points in the feature space, multiplies the distance by a random number between 0 and 1 and places the new data point at this new distance from one of the data points used for distance calculation. Note the number of nearest neighbors considered for data generation is also a hyperparameter and can be changed based on requirement.
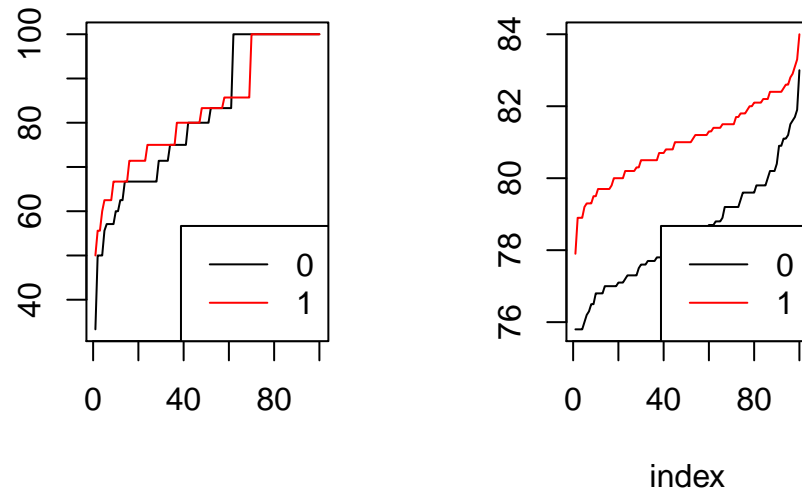
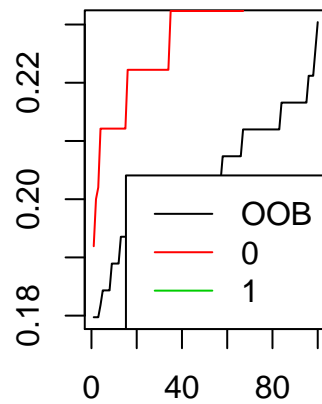**CV producers accuracy**   **Model producers accura**



Figure 13: RF Classifier cross-validation using original data

**CV oob error**          **Model oob error**
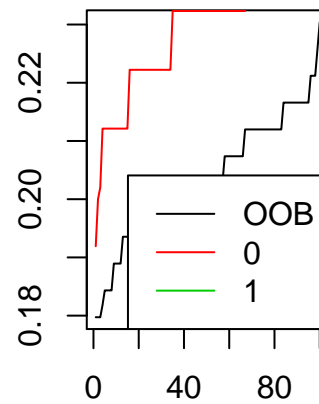


Figure 14: RF Classifier cross-validation OOB error using original data

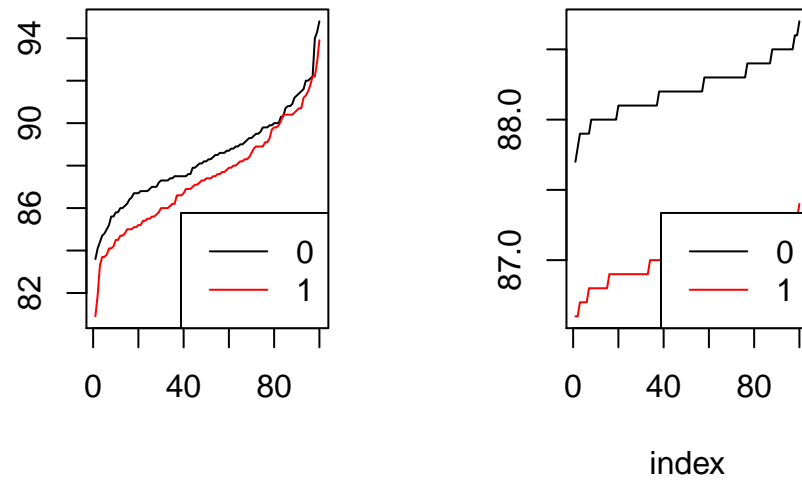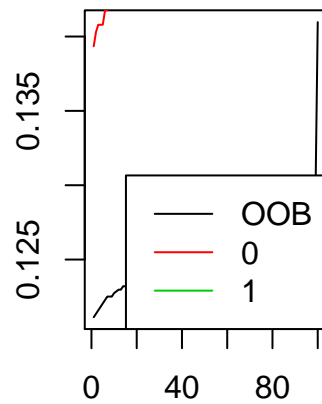**CV producers accuracy**  **Model producers accura**



Figure 15: RF Classifier cross-validation using synthesized data
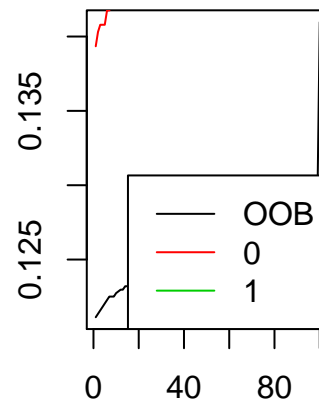
**CV oob error**  **Model oob error**



Figure 16: RF Classifier cross-validation OOB using synthesized data

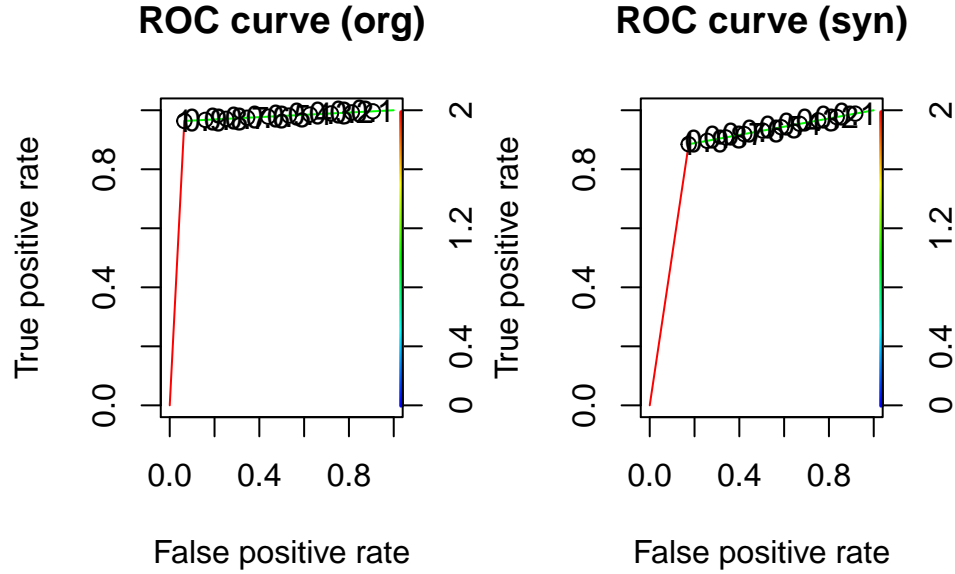**ROC curve (org)**          **ROC curve (syn)**

Figure 17: RF classifier ROC curve Original(left), Synthesized(right)

The chart shows accuracy and oob error graph of RF classifiers built on original data, synthesized data on top and bottom row respectively. It appears that the classifier built on original model shows higher OOB error despite the higher accuracy rate. This is due to classifier having `seen` the test data when building the model.

## 3.3 Support Vector Machines

Per our literature review, Assari et al found that Support Vector Machine (SVM) models provided better accuracy than Naive Bayes and Decision Tree approaches.

SVM models are a form of supervised learning used to create discriminative large-margin classifiers, defining a decision boundary between classes based on a hyperplane. This hyperplane maximizes the margin, or distance, between the nearest points of different labeled classes while prioritizing correct classification. The classifier can then be used to predictively categorize unlabeled examples.

SVM models perform well in high dimensions and can be used to produce linear as well as non-linear classification depending on the kernel function employed. However, they are prone to overfitting and can be hard to interpret, particularly when in higher dimensions.

### 3.3.1 Kernel Methods

Kernels are algorithms that operate in high-dimensional feature space without explicitly mapping data to the coordinates of that space. Mathematically, this 'kernel trick' means scalar solutions to dot products can be used without computationally laborious transformations.

Practically for SVM models, if we can map the feature space to higher dimensions in which the classes of cases become separable, then a linear classifier can be used to solve non-linear problems.

As kernel values depend on the inner products of feature vectors, it's best practice to scale variables to range of [0, 1] or [-1, 1], which can be accomplished with the preProcess argument of caret's train function.

Choice of kernel function has a significant impact on the outcome of an SVM model. While this provides flexibility it also create potential points of failure, and SVM models are very sensitive to selection of kernel parameters.

Hsu et al recommend commencing with the Guassian radial basis function (RBF) kernel, which can handle nonlinear relationships, has fewer hyperparameters, and prevents fewer numerical challenges. For this project, we evaluated both RBFs and linear kernels, setting them via the method argument of caret's train function.

### 3.3.2   Tuning

A hard margin constraint prevents the classifier from allowing the margin to overlap with values. This constraint can be relaxed, which reduces variance but adds bias. By manipulating the regularization, or penalty, parameter C we can decrease (soft margin) or increase (hard margin) the weight of values within the margin on overall error. To test the impact of different C values we employ a grid search using the expand.grid function.

The tuning parameter sigma also impacts model fit. A smaller sigma tends to yield a local classifier with less bias and more variance, while a larger sigma value tends to yield a general classifier with more bias and less variance.

### 3.3.3   Approach

SVM models can be computationally intensive, so we first build the model and evaluate performance based on the original dataset of 303 cases before attempting with a larger synthesized dataset. We harness different kernel functions and tuned parameters to optimize performance based on accuracy (number of correct predictions as percentage of total cases) and kappa statistics (comparing observed and expected accuracy based on random chance, or interrater reliability).

### 3.3.4   Models

For the first SVM model, the original dataset is preprocessed (centered and scaled) before applying an RBF kernel with cross-validation (10 folds, 5 repeats). caret identifies the optimal model based on ROC, which is highest for sigma = .033 and C = .25. This model yields an accuracy of .813 and a kappa statistic of .622.

We generate a second SVM model with same preprocessing, RBF kernel, and cross-validation as the first, tuning it via a grid search around the sigma and C values. This second model is optimal on ROC at a lower sigma value (.01) and higher regularization parameter (.4), displaying more variance and a harder decision boundary, but lower accuracy (.787) and kappa (.565).

As tuning did not improve accuracy, we apply the same approach and parameters of the first model (identical preprocessing, RBF kernel, and cross-validation) to the synthesized data to create a third SVM model. Based on ROC, caret selects the optimal model with highest sigma = .031 and C = 1. This model performs better than both preceding on accuracy (.84) and kappa (.672). Compared with the first model, sensitivity has declined (.765 vs. .735), meaning slightly more false negatives; but specificity has increased (.854 vs. .927), for fewer false positives.

For a contrasting approach, we also trial a linear classifier (in place of the radial used in the first three models) using the synthesized data to create a fourth SVM model. This model achieves better accuracy (.827) and kappa statistic (.65) than the first and second models but does not outperform the third.

In summary, an SVM model using a radial kernel atop the larger synthesized dataset achieved an accuracy of .84, with a stricter penalties on false positives than false negatives. Given the model is intended as a diagnostic health tool, the optimal balance between sensitivity and specificity could be further explored (i.e. chance of misdiagnosing a patient at higher risk vs. cost of needlessly testing a patient with lower risk).

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  X0  X1
```

```
##           X0 556  97
##           X1 151 711
##
##                 Accuracy : 0.8363
##                   95% CI : (0.8167, 0.8546)
##      No Information Rate : 0.5333
##      P-Value [Acc > NIR] : < 2.2e-16
##
##                    Kappa : 0.6696
##
##   Mcnemar's Test P-Value : 0.000764
##
##              Sensitivity : 0.7864
##              Specificity : 0.8800
##           Pos Pred Value : 0.8515
##           Neg Pred Value : 0.8248
##               Prevalence : 0.4667
##           Detection Rate : 0.3670
##     Detection Prevalence : 0.4310
##        Balanced Accuracy : 0.8332
##
##         'Positive' Class : X0
##
```

## 3.4   Naive Bayes

Naive Bayes classifier assumes that the presence (or absence) of a particular feature is unrelated to the presence (or absence) of any other feature. It considers all variables to independently contribute to the probability of heart disease. In spite of their naive design and apparently oversimplified assumptions, naive Bayes classifiers often work much better in many complex real world situations. Additionally, it requires a small amount of training data to estimate the parameters.

The most accurate model was build on synthetic data. The model looked at all categorical(factorized) variables. We removed numeric variables `age` and `sex` from the classifier to improve our model. Additionally, `chol` variable was converted into a categorical variable.

With the above parameters we are able to lift our accuracy to 88%, which is reasonably high. The Naive Bayes classifier is often hard to beat in terms of CPU and memory consumption, and in certain cases its performance can be very close to more complicated and slower techniques we used previously.

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction X0 X1
##         X0 28  0
##         X1  6 41
##
##                 Accuracy : 0.92
##                   95% CI : (0.834, 0.9701)
##      No Information Rate : 0.5467
##      P-Value [Acc > NIR] : 1.555e-12
##
##                    Kappa : 0.8361
##
##   Mcnemar's Test P-Value : 0.04123
```

```
##
##             Sensitivity : 0.8235
##             Specificity : 1.0000
##          Pos Pred Value : 1.0000
##          Neg Pred Value : 0.8723
##              Prevalence : 0.4533
##          Detection Rate : 0.3733
##    Detection Prevalence : 0.3733
##       Balanced Accuracy : 0.9118
##
##        'Positive' Class : X0
##
```

# 4 MODEL REVIEW AND SELECTION

## 4.1 Comparison of performance between models

[sensitivity, specificity, accuracy, others metrics?]

[Between different techniques and between original dataset and synthesize dataset for each technique]

## 4.2 Comparison of performance viz. other studies

# 5 CONCLUSIONS

# 6  APPENDIX

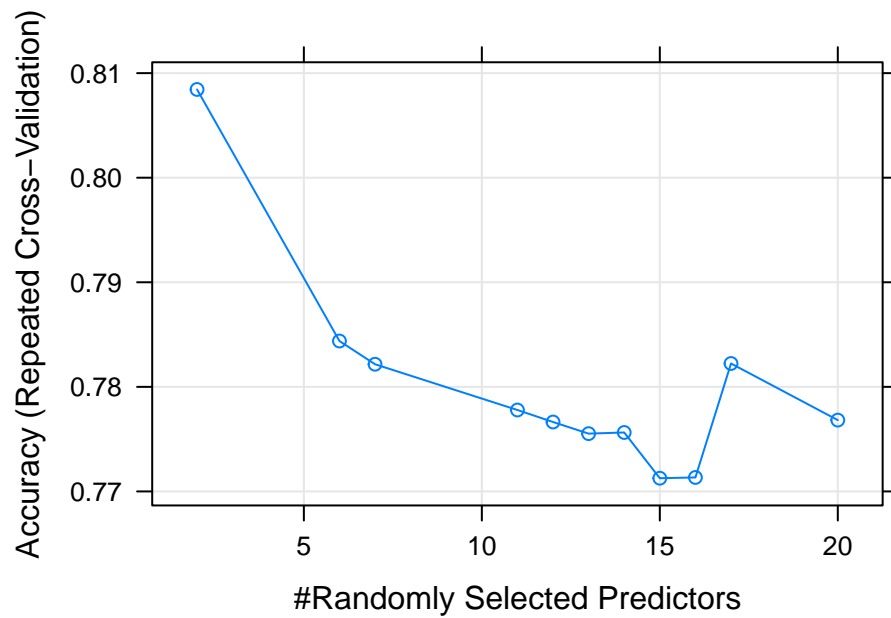## 6.1  Supplemental tables and figures



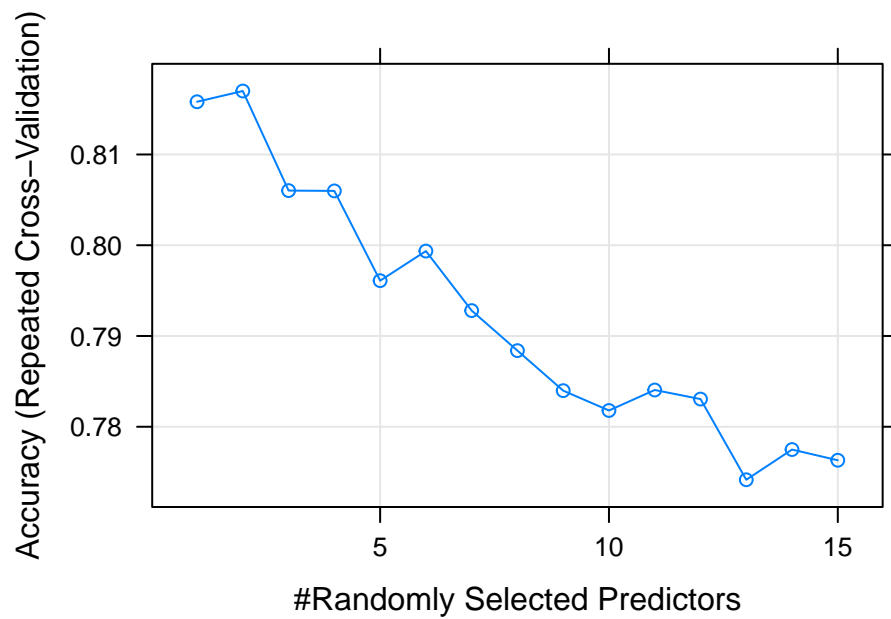Figure 18: Random Forest Hyperparameter Tuning - Random Search



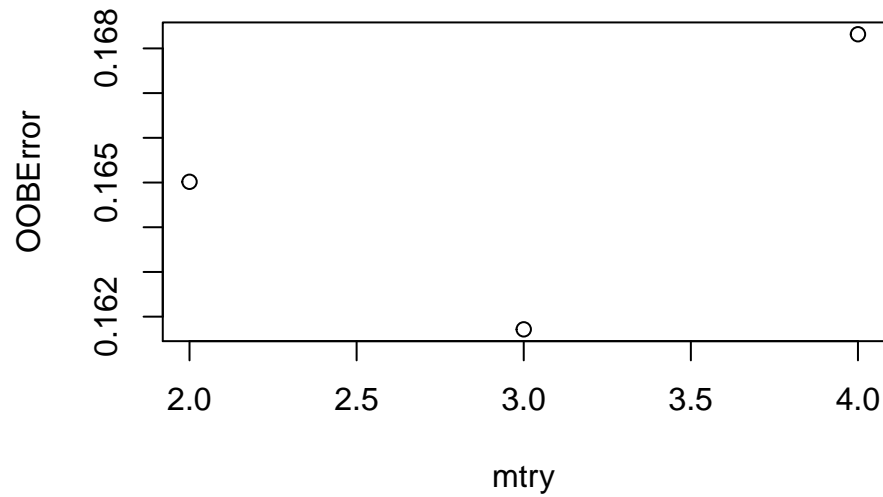Figure 19: Random Forest Hyperparameter Tuning - Grid Search

Figure 20: Random Forest Best mtry search

## 6.2 R statistical programming code

The appendix is available as script.R file in `projectFinal_heart` folder.

https://github.com/betsyrosalen/DATA_621_Business_Analyt_and_Data_Mining

Table 6: Previous Approaches Documented by Shouman

| AUTHOR | YEAR | TECHNIQUE | ACCURACY |
|--------|------|-----------|----------|
| Hall | 2000 | Na?ve Bayes | .832 |
| Hall | 2000 | K Nearest Neighbour | .821 |
| Hall | 2000 | Decision Tree | .753 |
| Yan, Zheng et al. | 2003 | Multilayer Perceptron | .636 |
| Herron | 2004 | Support Vector Machine | .836 |
| Herron | 2004 | J4.8 Decision Tree | .776 |
| Herron | 2004 | Support Vector Machine | .834 |
| Andreeva | 2006 | Na?ve Bayes | .786 |
| Andreeva | 2006 | Decision Tree | .757 |
| Andreeva | 2006 | Neural Network | .828 |
| Andreeva | 2006 | Sequential Minimal Optimization | .841 |
| Andreeva | 2006 | Kernel Density | .844 |
| Polat , Sahan et al. | 2007 | Fuzzy-AIRS-K-Nearest Neighbour | .870 |
| Palaniappan and Awang | 2007 | Na?ve Bayes | .950 |
| Palaniappan and Awang | 2007 | Decision Tree | .949 |
| Palaniappan and Awang | 2007 | Neural Network | .935 |
| De Beule, Maesa et al. | 2007 | Artificial Neural Network | .820 |
| Tantimongcolwat, Naenna et al. | 2008 | Direct Kernel Self-organizing Map | .804 |
| Tantimongcolwat, Naenna et al. | 2008 | Multilayer Perceptron | .745 |
| Hara and Ichimura | 2008 | Automatically Defined Groups | .678 |
| Hara and Ichimura | 2008 | Immune Multi-agent Neural Network | .823 |
| Sitar-Taut, Zdrenghea et al. | 2009 | Na?ve Bayes | .620 |
| Sitar-Taut, Zdrenghea et al. | 2009 | Decision Trees | .604 |
| Tu, Shin et al. | 2009 | Bagging Algorithm | .814 |
| Das, Turkoglu et al. | 2009 | Neural Network Ensembles | .890 |
| Rajkumar and Reena | 2010 | Na?ve Bayes | .523 |
| Rajkumar and Reena | 2010 | K Nearest Neighbour | .457 |
| Rajkumar and Reena | 2010 | Decision List | .520 |
| Srinivas, Rani et al. | 2010 | Na?ve Bayes | .841 |
| Srinivas, Rani et al. | 2010 | One Dependency Augmented Na?ve Bayes | .805 |
| Kangwanariyakul, Nantasenamat et al. | 2010 | Back-propagation Neural Network | .784 |
| Kangwanariyakul, Nantasenamat et al. | 2010 | Bayesian Neural Network | .784 |
| Kangwanariyakul, Nantasenamat et al. | 2010 | Probabilistic Neural Network | .706 |
| Kangwanariyakul, Nantasenamat et al. | 2010 | Polynomial Support Vector Machine | .706 |
| Kangwanariyakul, Nantasenamat et al. | 2010 | Radial Basis Support Vector Machine | .608 |
| Kangwanariyakul, Nantasenamat et al. | 2010 | Bayesian Neural Network | .784 |
| Kumari and Godara | 2011 | RIPPER | .811 |
| Kumari and Godara | 2011 | Decision Tree | .791 |
| Kumari and Godara | 2011 | Artificial Neural Network | ,801 |
| Kumari and Godara | 2011 | Support Vector Machine | .841 |
| Soni, Ansari et al. | 2011 | Weighted Associative Classifier | .578 |
| Soni, Ansari et al. | 2011 | Classification-Association | .583 |
| Soni, Ansari et al. | 2011 | Classification-Multiple ClassAssociation | .536 |
| Soni, Ansari et al. | 2011 | Classification-Predictive Association | .523 |
| Abdullah and Rajalaxmi | 2012 | Decision Tree | .507 |
| Abdullah and Rajalaxmi | 2012 | Random Forest | .633 |
| Rajeswari, Vaithiyanathan et al. | 2013 | Neural Network | .805 |
| Rajeswari, Vaithiyanathan et al. | 2013 | J4.8 Decision Tree | .779 |
| Rajeswari, Vaithiyanathan et al. | 2013 | Support Vector Machine | .842 |
| Rajeswari, Vaithiyanathan et al. | 2013 | Feature Selection with Neural Network | .845 |
| Rajeswari, Vaithiyanathan et al. | 2013 | Feature Selection with Decision Tree | .842 |
| Rajeswari, Vaithiyanathan et al. | 2013 | Feature Selection with Support Vector Machine | .875 |
| Rajeswari, Vaithiyanathan et al. | 2013 | Neural Network | .805 |
| Lakshmi, Krishna et al. | 2013 | Support Vector Machine | .781 |
| Lakshmi, Krishna et al. | 2013 | Decision Tree | .847 |
| Lakshmi, Krishna et al. | 2013 | K Nearest Neighbour | .840 |

Table 7: Average Performance by Technique

| TECHNIQUE | MEDIAN ACCURACY |
|---|---|
| Logistic Regression | .855 |
| Random Forest | .724 |
| Support Vector Machine | .809 |
| Naive Bayes | .819 |

| Sensitivity | Specificity | Precision | Recall | F1 |
|---|---|---|---|---|
| 0.935 | 0.964 | 0.956 | 0.935 | 0.945 |
| 0.826 | 0.885 | 0.857 | 0.826 | 0.841 |